

# Fast-rate and optimistic-rate error bounds for $\ell_1$ -regularized regression

Rina Foygel and Nathan Srebro

August 2, 2011

## Abstract

We consider the prediction error of linear regression with  $\ell_1$  regularization when the number of covariates  $p$  is large relative to the sample size  $n$ . When the model is  $k$ -sparse and well-specified, and restricted isometry or similar conditions hold, the excess squared-error in prediction can be bounded on the order of  $\frac{\sigma^2 k \log(p)}{n}$ , where  $\sigma^2$  is the noise variance. Although these conditions are close to necessary for accurate *recovery* of the true coefficient vector, it is possible to guarantee good predictive accuracy under much milder conditions, avoiding the restricted isometry condition, but only ensuring an excess error bound of order  $\frac{k \log(p)}{n} + \sigma \sqrt{\frac{k \log(p)}{n}}$ . Here we show that this is indeed the best bound possible (up to logarithmic factors) without introducing stronger assumptions similar to restricted isometry.

## 1 Introduction

We consider a random design linear regression problem with  $p$  covariates:

$$y = x^T \beta^* + z$$

where  $x \in \mathbb{R}^p$  are random covariates with covariance matrix  $\Sigma$ ,  $z$  is random noise with  $\mathbb{E}[z^2] = \sigma^2$ , and  $\beta^* \in \mathbb{R}^p$  are the regression coefficients. For simplicity we take the response to be normalized,  $\mathbb{E}[y^2] = 1$  (otherwise all results scale accordingly).

We consider the problem of minimizing the prediction error

$$\mathbb{E} \left[ (y - x^T \beta)^2 \right]$$

based on an i.i.d. sample  $(x^{(1)}, y^{(1)}) \dots, (x^{(n)}, y^{(n)})$  using  $\ell_1$ -regularized regression:

$$\hat{\beta}^B \doteq \arg \min_{\|\beta\|_1 \leq B} \sum_i \left( y^{(i)} - x^{(i)T} \beta \right)^2.$$

Note that up to some unknown and data-dependent correspondence between  $B$  and  $\lambda$ , this is the same as

$$\hat{\beta}_\lambda \doteq \arg \min_{\beta} \sum_i \left( y^{(i)} - x^{(i)T} \beta \right)^2 + \lambda \|\beta\|_1,$$

also known as Lasso regression [Tibshirani, 1996].

Suppose that the covariates are 1-bounded, and that  $\max_i |y^{(i)}| \leq \mathbf{O}(\log(np))$  (for instance, this is true with high probability in the Gaussian setting). Then, by Srebro et al. [2010], with high probability over the sample, for any fixed  $\beta^*$  with  $\|\beta^*\|_1 \leq B$ , excess squared-error under  $\ell_1$ -regularized regression is bounded as

$$\mathbb{E} \left[ (y - x^T \hat{\beta}^B)^2 \right] - \mathbb{E} \left[ (y - x^T \beta^*)^2 \right] = \mathbf{O} \left( \frac{(1+B)^2 \log(p)}{n/\log^3(n)} + \sqrt{\frac{(1+B)^2 \log(p)}{n/\log^3(n)}} \cdot \mathbb{E} \left[ (y - x^T \beta^*)^2 \right] \right). \quad (1)$$

This result does not require any conditions on the correlation between the covariates, or on the nature of the “noise”  $y - x^T \beta^*$ , aside from the mild bound on  $\max_i |y^{(i)}|$ . In particular, this noise is not required to be independent from  $x$ . We believe also that this result would hold for subgaussian  $x$ ’s (rather than our current stronger assumption that the  $x$ ’s are 1-bounded).

We can apply this result to the sparse regression setting, with some mild additional assumptions. Suppose that we are interested in comparing to a sparse predictor on an unknown support  $J^* \subset [p]$ , with  $|J^*| \leq k$ . We now place a lower-bound eigenvalue assumption on this support  $J^*$  only:

$$\lambda_{\min}(\mathbb{E}[x_{J^*} x_{J^*}^T]) \geq \lambda_1 > 0, \quad (2)$$

where  $x_{J^*} = (x_j : j \in J^*)$  is the random vector consisting of those covariates  $x_j$  for which  $\beta_j^*$  is nonzero. This assumption is strictly weaker than the restricted isometry property (RIP) conditions in the compressed sensing literature, which require an upper-bound assumption as well, and require the eigenvalue bounds to hold for all sets  $J \subset [p]$  of bounded size, in addition to the true support  $J^*$ .

We fix the scale of the problem by assuming  $\mathbb{E}[y^2] = 1$ . Now consider a predictor  $\beta^*$  with support in  $S^*$ , which is better than the zero predictor — that is,  $\mathbb{E}[(y - x^T \beta^*)^2] \leq \mathbb{E}[(y - x^T \mathbf{0}_p)^2] = \mathbb{E}[y^2] = 1$ . We now show that  $\|\beta^*\|_1 = \mathbf{O}\left(\sqrt{k\lambda_1^{-1}}\right)$ . We first bound  $\|\beta^*\|_2^2$ , by observing that

$$\begin{aligned} \|\beta^*\|_2^2 \cdot \lambda_1 &\leq (\beta^*)^T \mathbb{E}[x x^T] \beta^* = \mathbb{E}[(x^T \beta^*)^2] = \mathbb{E}[(y - x^T \beta^*)^2] - 2\mathbb{E}[y \cdot (y - x^T \beta^*)] + \mathbb{E}[y^2] \\ &\leq 2\mathbb{E}[(y - x^T \beta^*)^2] + 2\mathbb{E}[y^2] \leq 4. \end{aligned}$$

We then have

$$\|\beta^*\|_1 \leq \sqrt{k} \|\beta^*\|_2 \leq \sqrt{k \cdot 4\lambda_1^{-1}} = \mathbf{O}\left(\sqrt{k\lambda_1^{-1}}\right).$$

Therefore, with high probability,

$$\mathbb{E} \left[ (y - x^T \hat{\beta}^B)^2 \right] - \mathbb{E} \left[ (y - x^T \beta^*)^2 \right] = \mathbf{O} \left( \frac{k \log(p)}{\lambda_1 n / \log^3(n)} + \sqrt{\frac{k \log(p)}{\lambda_1 n / \log^3(n)}} \cdot \mathbb{E} \left[ (y - x^T \beta^*)^2 \right] \right), \quad (3)$$

under the assumption that the  $x$ ’s are 1-bounded and  $\max_i |y^{(i)}| \leq \mathbf{O}(np)$ . Therefore, to guarantee a bound of  $\epsilon$  on the excess prediction error, the required sample complexity is

$$n = \Theta \left( \frac{k \log(p)}{\lambda_1 \epsilon} \cdot \frac{\sigma^2 + \epsilon}{\epsilon} \cdot \log^3(k/\lambda_1 \epsilon) \right), \quad (4)$$

where  $\sigma^2 = \mathbb{E}[(y - x^T \beta^*)^2]$  is the magnitude of the noise. This sample complexity follows an “optimistic rate”: in the noisy setting, if we would like to ensure a bound  $\epsilon$  on excess error which is small relative to  $\sigma^2$ ,

then the required sample complexity is then  $n = \Theta(\epsilon^{-2})$ , but on the other hand, in the noiseless setting (i.e. when  $y = x^T \beta^*$ ), or if the bound on excess error  $\epsilon$  is not much smaller than  $\sigma^2$ , then we require only  $n = \Theta(\epsilon^{-1})$ . We emphasize that this result does not assume that the linear model is a true model or require independent noise.

In contrast, results on sparse vector recovery from the compressed sensing framework [Candes and Tao, 2005, Bickel et al., 2009, Koltchinskii, 2009, Cai et al., 2009] provide stronger guarantees in a similar setting, using either  $\ell_1$ -regularized regression or the Dantzig selector, given by

$$\hat{\beta}_\lambda^{DS} = \arg \min \max_i \left| y^{(i)} - x^{(i)T} \beta \right| + \lambda \|\beta\|_1 .$$

These stronger results require several additional specialized assumptions, including the requirement that the noise must be independent from the signal. Existing results are stated either in the deterministic or random covariates setting, but can in general be translated to a random Gaussian setting. We restrict our attention to  $\ell_1$ -regularized regression when the covariates are i.i.d. multivariate Gaussian with zero mean:  $x^{(i)} \stackrel{i.i.d.}{\sim} N(0, \Sigma)$ . We now summarize this setting (with some simplifications), and compare it to the optimistic-rate results discussed above.

- **Well-specified model with independent subgaussian noise:** Response  $y^{(i)}$  is given by  $y^{(i)} = x^{(i)T} \beta^* + \sigma z^{(i)}$  for a true predictor  $\beta^*$  satisfying  $\|\beta^*\|_1 \leq B$ , and  $z^{(i)}$  is a subgaussian or subexponential noise term with unit variance, and is independent from  $x^{(i)}$ .

The main additional requirement here is that noise  $z$  is independent of  $x$ . This in particular implies that  $\beta^*$  is the optimal regressor. Note that in order to obtain the optimistic-rate guarantee (3), no such assumption is necessary, and  $\beta^*$  can be a non-optimal regressor chosen for its sparsity or eigenvalue properties.

- **Sparsity:**  $\beta^*$  is  $k$ -sparse, meaning that it has (at most)  $k$  non-zero entries.  
To obtain the optimistic-rate guarantee as stated originally in (1), we can relax this requirement and only assume that  $\beta^*$  has low  $\ell_1$ -norm.
- **Restricted eigenvalues:** There exists a  $\kappa \doteq \kappa(k, 3) > 0$ , such that for any  $J \subset [p]$  with  $|J| \leq k$ , for any nonzero  $\beta \in \mathbb{R}^p$  with  $\|\beta_{-J}\|_1 \leq 3\|\beta_J\|_1$ ,

$$\beta^T \Sigma \beta \geq \kappa \|\beta_J\|_2^2 . \quad (5)$$

This restricted eigenvalue condition is implied by a stronger condition:

**Restricted isometry:** Suppose that  $\delta_{2k} + 3\theta_{k,2k} < 1$ , where  $v^T \Sigma v \in (1 \pm \delta_{2k}) \|v\|_2^2$  for all  $2k$ -sparse vectors  $v$ , and  $|v^T \Sigma w| \leq \theta_{k,2k} \|v\|_2 \|w\|_2$  for all  $k$ -sparse  $v$  and  $2k$ -sparse  $w$  with disjoint supports. Then  $\kappa \doteq \sqrt{1 - \delta_{2k}} \left(1 - \frac{3\theta_{k,2k}}{1 - \delta_{2k}}\right)$  satisfies the restricted eigenvalue condition above.

To obtain the optimistic-rate guarantee (3) under the sparsity assumption, we required an eigenvalue condition (2) on  $\Sigma_{\text{Support}(\beta^*)}$  only, which is strictly weaker than the restricted eigenvalue and restricted isometry assumptions.

Under these assumptions, with  $\kappa$  defined as in (5), the following guarantees hold with high probability, by Theorem 7.2 of Bickel et al. [2009]:

$$\begin{aligned} \text{Sparse and accurate estimation of } \beta^*: \quad & \left\| \hat{\beta}^B - \beta^* \right\|_1 = \mathcal{O} \left( \frac{\sigma k}{\kappa^2} \cdot \sqrt{\frac{\log(p)}{n}} \right), \text{ and } \|\hat{\beta}^B\|_0 = \mathcal{O}(k) . \\ \text{Bounded excess prediction error: } & \mathbb{E} \left[ \left( y - x^T \hat{\beta}^B \right)^2 \right] = \mathbb{E} \left[ \left( y - x^T \beta^* \right)^2 \right] + \mathcal{O} \left( \frac{\sigma^2 k \log(p)}{\kappa^2 n} \right) . \end{aligned} \quad (6)$$

This corresponds to a sample complexity of

$$n = \Theta \left( \frac{\sigma^2 k \log(p)}{\kappa^2 \epsilon} \right), \quad (7)$$

to ensure an excess error bound of  $\epsilon$ . It is crucial to note that the error bound (and the sample complexity) scales with the magnitude of the noise,  $\sigma^2$ , rather than to the (unit) magnitude of the signal. In particular, in a noiseless setting, the results above guarantee a zero-error reconstruction of  $\beta^*$ , in contrast to the “optimistic rate” result (3) where no such guarantee is given. Furthermore, in the noisy setting, the compressed sensing guarantees give a “fast rate” result, since the sample complexity scales with  $\frac{1}{\epsilon}$  rather than with  $\frac{1}{\epsilon^2}$ .

In this compressed sensing framework, the guarantees on predictive error follow from a stronger guarantee on the accurate recovery of  $\beta^*$ , and in particular, the recovery of the true support of  $\beta^*$ . In order for this to be possible, it is of course necessary to be able to distinguish between pairs or small sets of covariates. In particular, some sort of restricted isometry assumption is clearly necessary for bounding error in recovering  $\beta^*$  (otherwise, the “best”  $\beta^*$  might not be unique). However, if the goal is merely low error in prediction — that is, we would like accuracy in calculating  $x^T \beta^*$ , rather than in recovering  $\beta^*$  — then perhaps this assumption could be weakened. For example, if a covariate is duplicated in the model, then it will not be possible to distinguish between the two when attempting to recover the true support; however, adding duplicated covariates to a model will have no effect on the problem of prediction.

More generally, we are interested in whether the properties that are necessary for the (unique) recovery of  $\beta^*$ , are also necessary to obtain strong bounds on excess prediction error, and in the role of the assumptions that separate the “optimistic rate”, unit-scale error bounds of Srebro et al. [2010] from the “fast rate” error bounds in the compressed sensing literature, which scale with the magnitude of the noise. Below, we show that, if we remove either the sparsity assumption (while still assuming that  $\beta^*$  has low  $\ell_1$ -norm) or the restricted isometry assumption from the compressed sensing framework described above, then up to logarithmic factors, the “optimistic rate” bound on excess prediction error, given in (3), is the best possible bound. In particular, this implies that, even in the noiseless setting, we cannot achieve zero error in prediction, without stronger assumptions.

## 2 Results

First, we ask whether we can relax the assumption of a sparse true coefficient vector to an assumption on its  $\ell_1$ -norm, but still guarantee a fast-rate bound on excess error. Specifically, we consider the question of bounding excess prediction error, in the well-specified Gaussian setting where the restricted eigenvalue assumption holds, assuming only an  $\ell_1$ -norm bound on the true vector of coefficients.

Our first result shows that, up to logarithmic factors, the optimistic-rate error bound (3) is the best possible rate under these conditions. For simplicity, we will consider the case of completely independent covariates,  $x \sim N(0, \mathbf{I}_p)$ . In particular, this ensures that the restricted eigenvalue assumption is satisfied. To place the problem on a unit scale (or rather, to bound the scale away from zero and away from infinity), we consider only true coefficient vectors  $\beta^*$  satisfying

$$\frac{1}{2} \leq \mathbb{E} \left[ (x^T \beta^*)^2 \right]^{1/2} = \|\beta^*\|_2 \leq \|\beta^*\|_1 \leq 1.$$

**Theorem 1.** *Fix any  $n \geq 30$ ,  $p \geq 3n$ , and  $\sigma \geq 0$ . Then there exists a  $\beta^* \in \mathbb{R}^p$  with  $\frac{1}{2} \leq \|\beta^*\|_2 \leq \|\beta^*\|_1 \leq 1$ , such that for any sample, for all  $B \geq 0$ ,*

$$\mathbb{E} \left[ \left( y - x^T \hat{\beta}^B \right)^2 \right] \geq \sigma^2 + \frac{1}{32n \log^2(3n)}.$$

Additionally, if  $100 \leq \sqrt{n}/\sigma \leq p$ , then with probability at least  $\frac{1}{2}$  over the sample, for all  $B \geq 0$ ,

$$\mathbb{E} \left[ \left( y - x^T \hat{\beta}^B \right)^2 \right] \geq \sigma^2 + \frac{\sigma}{102400 \sqrt{n} \log^2 (\max \{3n, \lceil \sqrt{n}/\sigma \rceil\})}.$$

Here  $\hat{\beta}^B = \arg \min_{\|\beta\|_1 \leq B} \sum_i (y^{(i)} - x^{(i)T} \beta)^2$ , where  $(x^{(i)}, y^{(i)})$  are i.i.d. samples from the multivariate Gaussian distribution defined by drawing  $x^{(i)} \sim N(0, \mathbf{I}_p)$  and  $y^{(i)} \sim N(x^{(i)T} \beta^*, \sigma^2)$ . The expectations are taken over a new sample  $(x, y)$  drawn from the same distribution, independently of the training set  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$ . (For each  $B \geq 0$ , if  $\hat{\beta}^B$  is not unique, then we show that the inequalities hold for some choice of  $\hat{\beta}^B$ .)

Next, we ask whether the restricted eigenvalue (or restricted isometry) assumption is necessary for a fast-rate bound on excess error, in the well-specified Gaussian setting where the sparsity assumption holds.

Our second result shows that, up to logarithmic factors, the optimistic-rate error bound (3) is the best possible rate under these conditions. For simplicity, we restrict our attention to 2-sparse true coefficient vectors. We also only consider covariance matrices  $\Sigma$  such that  $\Sigma_{J^*} = \mathbf{I}_{J^*}$ , where  $J^* = \text{Support}(\beta^*)$ . That is, ensuring the restricted isometry property on the true support only, is not sufficient for a fast-rate bound on excess error.

To avoid issues of scaling, we restrict our attention to covariance matrices  $\Sigma$  with  $\|\Sigma\|_{sp} \leq 2$ , and to true coefficient vectors  $\beta^*$  satisfying

$$\frac{1}{2} \leq \mathbb{E} \left[ (x^T \beta^*)^2 \right]^{1/2} = \sqrt{\beta^{*T} \Sigma \beta^*} = \|\beta^*\|_2 \leq \|\beta^*\|_1 \leq 1,$$

where we make use of the fact that  $\Sigma_{J^*} = \mathbf{I}_{J^*}$  to obtain the second equality.

**Theorem 2.** Fix any  $n \geq 30$ ,  $p \geq 3n$ , and  $\sigma \geq 0$ . Then there exists a 2-sparse  $\beta^* \in \mathbb{R}^p$  with  $\frac{1}{2} \leq \|\beta^*\|_2 \leq \|\beta^*\|_1 \leq 1$ , and a positive semi-definite  $\Sigma \in \mathbb{R}^{p \times p}$  with  $\|\Sigma\|_{sp} \leq 2$  and  $\Sigma_{\text{Support}(\beta^*)} = \mathbf{I}_{\text{Support}(\beta^*)}$ , such that for any sample, for all  $B \geq 0$ ,

$$\mathbb{E} \left[ \left( y - x^T \hat{\beta}^B \right)^2 \right] \geq \sigma^2 + \frac{1}{288n \log^2(3n)}.$$

Additionally, if  $100 \leq \sqrt{n}/\sigma \leq p - 3$ , then with probability at least  $\frac{1}{2}$  over the sample, for all  $B \geq 0$ ,

$$\mathbb{E} \left[ \left( y - x^T \hat{\beta}^B \right)^2 \right] \geq \sigma^2 + \frac{\sigma}{409600 \sqrt{n} \log^2 (\max \{3n, \lceil \sqrt{n}/\sigma \rceil\})}.$$

Here  $\hat{\beta}^B = \arg \min_{\|\beta\|_1 \leq B} \sum_i (y^{(i)} - x^{(i)T} \beta)^2$ , where  $(x^{(i)}, y^{(i)})$  are i.i.d. samples from the multivariate Gaussian distribution defined by drawing  $x^{(i)} \sim N(0, \Sigma)$  and  $y^{(i)} \sim N(x^{(i)T} \beta^*, \sigma^2)$ . The expectations are taken over a new sample  $(x, y)$  drawn from the same distribution, independently of the training set  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$ . (For each  $B \geq 0$ , if  $\hat{\beta}^B$  is not unique, then we show that the inequalities hold for some choice of  $\hat{\beta}^B$ .)

In particular, Theorem 2 shows that without placing any assumptions on the covariates outside of  $\text{Support}(\beta^*)$ , we cannot guarantee a bound on excess error that is better than the optimistic rate obtained by Srebro et al. [2010] from concentration bounds, up to logarithmic factors.

### 3 Proofs

We begin by defining a class of predictors that are optimal with respect to the squared-error loss and the  $\ell_1$ -norm regularizer:

**Definition 1.** Given  $y^{(i)} \in \mathbb{R}$  and  $x^{(i)} \in \mathbb{R}^p$  for  $i = 1, \dots, n$ , a predictor  $\tilde{\beta} \in \mathbb{R}^p$  is Pareto-optimal (with respect to empirical squared-error and  $\ell_1$ -norm) if it satisfies

$$\sum_i \left( y^{(i)} - x^{(i)T} \beta \right)^2 \leq \sum_i \left( y^{(i)} - x^{(i)T} \tilde{\beta} \right)^2 \Rightarrow \|\beta\|_1 \geq \|\tilde{\beta}\|_1 ,$$

that is, if we cannot improve its empirical squared error without increasing its  $\ell_1$ -norm, and vice versa.

The following lemma states a well-known property of  $\ell_1$ -regularized regression; we include a proof for completeness.

**Lemma 1.** For any  $y^{(1)}, \dots, y^{(n)} \in \mathbb{R}$  and  $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^p$ , for any  $B \geq 0$ , the class

$$\mathcal{B}_B \doteq \arg \min_{\|\beta\|_1 \leq B} \sum_i \left( y^{(i)} - x^{(i)T} \beta \right)^2$$

must contain a predictor  $\hat{\beta}^B$  that is Pareto-optimal and satisfies  $\|\hat{\beta}^B\|_0 \leq n$ .

*Proof.* Let  $\text{Err}_B = \inf_{\|\beta\|_1 \leq B} \sum_i \left( y^{(i)} - x^{(i)T} \beta \right)^2$ . Since  $\{\|\beta\|_1 \leq B\}$  is a compact set, this infimum is attained by some  $\beta$  with  $\|\beta\|_1 \leq B$ . Now define

$$B' = \inf \left\{ \|\beta\|_1 : \sum_i \left( y^{(i)} - x^{(i)T} \beta \right)^2 \leq \text{Err}_B \right\} \leq B .$$

Again, by compactness, this infimum is attained by some  $\tilde{\beta}$ . We then see that  $\tilde{\beta}$  is Pareto-optimal by its construction. Finally, by Theorem 3 of Rosset et al. [2004], there exists a  $\hat{\beta}^B \in \mathbb{R}^p$  such that  $\|\hat{\beta}^B\|_0 \leq n$ ,  $\|\hat{\beta}^B\|_1 \leq \|\tilde{\beta}\|_1$ , and  $X\hat{\beta}^B = X\tilde{\beta}$ . This is sufficient.  $\square$

Next we state two additional lemmas, proved in the next section.

**Lemma 2.** Fix  $n$  and  $p$  with  $n \geq 30$  and  $p \geq 3n$ . Let  $x^{(i)} \stackrel{i.i.d.}{\sim} N(0, \Sigma)$  for some  $\Sigma \in \mathbb{R}^{p \times p}$ , and let  $\beta^* \in \mathbb{R}^p$  be fixed. Then with probability at least  $1 - 2e^{-n \log(p)}$ , for all  $J \subset [p]$  with  $|J| = n$ ,

$$\left\| X_J^T X (\tilde{\beta} - \beta^*) \right\|_2 \leq \|\Sigma\|_{sp} \cdot 16\sqrt{2} \cdot n \log(p) \cdot \sqrt{(\tilde{\beta} - \beta^*)^T \Sigma (\tilde{\beta} - \beta^*)} \text{ for all } \tilde{\beta} \in \mathbb{R}^p \text{ with } \tilde{\beta}_{-J} = 0 , \quad (8)$$

where the matrix  $X$  has entries  $X_{ij} = x_j^{(i)}$ , and  $X_J$  consists of the columns of  $X$  indexed by  $j \in J$ .

**Lemma 3.** Let  $x^{(i)} \stackrel{i.i.d.}{\sim} N(0, \Sigma)$  for some  $\Sigma \in \mathbb{R}^{p \times p}$ , and let  $z \in \mathbb{R}^n$  be fixed, with  $\|z\|_2^2 \geq 0.5n$ . Assume  $\sqrt{n}/\sigma \geq 100$ . Then with probability at least  $1 - e^{-0.015\sigma^{-1}\sqrt{n}}$ , for all  $J_1 \subset [\lceil \sqrt{n}/\sigma \rceil]$  with  $|J_1| \geq \frac{\sqrt{n}}{2\sigma}$ ,

$$\left\| \text{Proj}_{\mathbf{1}_{J_1}}^\perp (X_{J_1}^T z) \right\|_2^2 \geq \frac{\lambda_{\min}^2(\Sigma_{J_1}) n^{3/2}}{200\sigma} , \quad (9)$$

where the matrix  $X$  has entries  $X_{ij} = x_j^{(i)}$ , and  $X_{J_1}$  consists of the columns of  $X$  indexed by  $j \in J_1$ .

We now prove the theorems.

### 3.1 Proof of Theorem 1

*Theorem 1.* Fix any  $n \geq 30$ ,  $p \geq 3n$ , and  $\sigma \geq 0$ . Then there exists a  $\beta^* \in \mathbb{R}^p$  with  $\frac{1}{2} \leq \|\beta^*\|_2 \leq \|\beta^*\|_1 \leq 1$ , such that for any sample, for all  $B \geq 0$ ,

$$\mathbb{E} \left[ \left( y - x^T \hat{\beta}^B \right)^2 \right] \geq \sigma^2 + \frac{1}{32n \log^2(3n)}. \quad (10)$$

Additionally, if  $100 \leq \sqrt{n}/\sigma \leq p$ , then with probability at least  $\frac{1}{2}$  over the sample, for all  $B \geq 0$ ,

$$\mathbb{E} \left[ \left( y - x^T \hat{\beta}^B \right)^2 \right] \geq \sigma^2 + \frac{\sigma}{102400 \sqrt{n} \log^2(\max\{3n, \lceil \sqrt{n}/\sigma \rceil\})}. \quad (11)$$

Here  $\hat{\beta}^B = \arg \min_{\|\beta\|_1 \leq B} \sum_i (y^{(i)} - x^{(i)T} \beta)^2$ , where  $(x^{(i)}, y^{(i)})$  are i.i.d. samples from the multivariate Gaussian distribution defined by drawing  $x^{(i)} \sim N(0, \mathbf{I}_p)$  and  $y^{(i)} \sim N(x^{(i)T} \beta^*, \sigma^2)$ . The expectations are taken over a new sample  $(x, y)$  drawn from the same distribution, independently of the training set  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$ . (For each  $B \geq 0$ , if  $\hat{\beta}^B$  is not unique, then we show that the inequalities hold for some choice of  $\hat{\beta}^B$ .)

*Proof.* Let  $\beta^*$  be

$$\beta_j^* = \frac{1}{j \cdot 4 \log p}, \quad j = 1, \dots, p-1; \quad \beta_p^* = \frac{1}{2}.$$

Note that  $\|\beta^*\|_1 \leq 1$  and  $\|\beta^*\|_2^2 \geq \frac{1}{4}$ , and so the resulting distribution satisfies the desired assumptions.

By Lemma 1, for any  $B \geq 0$ , the set  $\arg \min_{\|\beta\|_1 \leq B} \sum_i (y^{(i)} - x^{(i)T} \beta)^2$  must include a Pareto-optimal vector  $\hat{\beta}^B$  with  $\|\hat{\beta}^B\|_0 \leq n$ . Therefore, it is sufficient to show that bounds (10) and (11) hold for all Pareto-optimal vectors  $\hat{\beta}$  with  $\|\hat{\beta}\|_0 \leq n$ . We now prove these two bounds separately.

**Proof of (10).** For any  $\hat{\beta}$  with  $\|\hat{\beta}\|_0 \leq n$ , we have

$$\begin{aligned} \left\| \hat{\beta} - \beta^* \right\|_2^2 &\geq \sum_{j=1}^p \left( \hat{\beta}_j - \frac{1}{j \cdot 4 \log p} \right)^2 \geq \sum_{j: \hat{\beta}_j = 0} \left( \frac{1}{j \cdot 4 \log p} \right)^2 \geq \sum_{j=n+1}^p \left( \frac{1}{j \cdot 4 \log p} \right)^2 \\ &\geq \frac{1}{16 \log^2(p)} \int_{x=n+1}^p \frac{1}{x^2} dx = \frac{1}{16 \log^2(p)} \left( \frac{1}{n+1} - \frac{1}{p} \right) \geq \frac{1}{32n \log^2(p)}. \end{aligned}$$

This proves the claim when  $p = 3n$ . However, the claim is immediately true for any larger value of  $p$ , since we may add in an arbitrary number of zero covariates (and assign zero coefficients to these covariates), without affecting the results.

**Proof of (11).** By Lemma 1 of Laurent and Massart [2000], with probability at least  $1 - e^{-0.0625n} \geq 0.75$ ,  $\|z\|_2^2 \sim \chi_n^2 \geq 0.5n$ . For the remainder of the proof, we treat  $z \in \mathbb{R}^n$  as a fixed vector, and assume  $\|z\|_2^2 \geq 0.5n$ .

Assume that (8) holds for all  $J \subset [p]$  with  $|J| = n$ , and (9) holds for all  $J_1 \subset [\lceil \sqrt{n}/\sigma \rceil]$  with  $|J_1| \geq \sqrt{n}/2\sigma$ . (By Lemmas 2 and 3, this is true with probability at least  $1 - 2e^{-n \log(p)} - e^{-0.015\sigma^{-1}\sqrt{n}} \geq 0.75$ .) Now choose any Pareto-optimal  $\hat{\beta}$  with  $\|\hat{\beta}\|_0 \leq n$ .

Suppose that  $\|\hat{\beta} - \beta^*\|_2^2 < \frac{\sigma}{102400\sqrt{n}\log^2(p)}$ . First, we show that

$$\left| \left\{ j \in [\lceil \sqrt{n}/\sigma \rceil] : \hat{\beta}_j > 0 \right\} \right| \geq \frac{\sqrt{n}}{2\sigma}.$$

Suppose not. Then

$$\begin{aligned} \|\hat{\beta} - \beta^*\|_2^2 &\geq \sum_{j \in [\lceil \sqrt{n}/\sigma \rceil] : \hat{\beta}_j \leq 0} (\beta_j^*)^2 = \frac{1}{16\log^2(p)} \sum_{j \in [\lceil \sqrt{n}/\sigma \rceil] : \hat{\beta}_j \leq 0} \frac{1}{j^2} \geq \frac{1}{16\log^2(p)} \sum_{j=\lceil \sqrt{n}/2\sigma \rceil}^{\lceil \sqrt{n}/\sigma \rceil} \frac{1}{j^2} \\ &\geq \frac{1}{16\log^2(p)} \int_{x=\lceil \sqrt{n}/2\sigma \rceil}^{2\lceil \sqrt{n}/2\sigma \rceil} \frac{1}{x^2} dx = \frac{1}{16\log^2(p)} \left( \frac{1}{\lceil \sqrt{n}/2\sigma \rceil} - \frac{1}{2\lceil \sqrt{n}/2\sigma \rceil} \right) = \frac{1}{16\log^2(p) \cdot 2\lceil \sqrt{n}/2\sigma \rceil} \geq \frac{\sigma}{32\sqrt{n}\log^2(p)}. \end{aligned}$$

This is a contradiction.

Now define  $J_1 = \{j \in [\lceil \sqrt{n}/\sigma \rceil] : \hat{\beta}_j > 0\}$ , and fix any  $J \supset \text{Support}(\hat{\beta})$  with  $|J| = n$ . Since  $\hat{\beta}$  is Pareto-optimal with positive entries  $\hat{\beta}_j$  for all  $j \in J_1$ , we have

$$\frac{\partial}{\partial(\beta_{J_1})} \|\hat{\beta}\|_1 = \mathbf{1}_{J_1}.$$

Therefore, by the theory of Lagrange multipliers, we must have  $X_{J_1}^T y - X_{J_1}^T X \hat{\beta} = C \cdot \mathbf{1}_{J_1}$ , for some  $C \in \mathbb{R}$ . We then have

$$X_{J_1}^T X(\hat{\beta} - \beta^*) = \sigma \cdot X_{J_1}^T z - C \cdot \mathbf{1}_{J_1}. \quad (12)$$

By (8), the norm of the left-hand side of (12) can be bounded from above as

$$\left\| X_{J_1}^T X(\hat{\beta} - \beta^*) \right\|_2 \leq \left\| X_{J_1}^T X(\hat{\beta} - \beta^*) \right\|_2 \leq \|\Sigma\|_{sp} \cdot 16\sqrt{2} \cdot n \log(p) \cdot \sqrt{(\hat{\beta} - \beta^*)^T \Sigma (\hat{\beta} - \beta^*)}.$$

By (9), the norm of the right-hand side of (12) can be bounded from below as

$$\left\| \sigma \cdot X_{J_1}^T z - C \cdot \mathbf{1}_{J_1} \right\|_2 \geq \sigma \cdot \left\| \text{Proj}_{\mathbf{1}_{J_1}^\perp} X_{J_1}^T z \right\|_2 \geq \sigma \sqrt{\frac{\lambda_{\min}^2(\Sigma_{J_1}) n^{3/2}}{200\sigma}}.$$

Therefore, returning to (12), we have

$$\begin{aligned} \|\Sigma\|_{sp} \cdot 16\sqrt{2} \cdot n \log(p) \cdot \sqrt{(\hat{\beta} - \beta^*)^T \Sigma (\hat{\beta} - \beta^*)} &\geq \left\| X_{J_1}^T X(\hat{\beta} - \beta^*) \right\|_2 \\ &= \left\| \sigma \cdot X_{J_1}^T z - C \cdot \mathbf{1}_{J_1} \right\|_2 \geq \sigma \sqrt{\frac{\lambda_{\min}^2(\Sigma_{J_1}) n^{3/2}}{200\sigma}}. \end{aligned}$$

Therefore,

$$(\hat{\beta} - \beta^*)^T \Sigma (\hat{\beta} - \beta^*) \geq \frac{\sigma \cdot \lambda_{\min}^2(\Sigma_{J_1})}{102400 \|\Sigma\|_{sp}^2 \cdot \sqrt{n} \log^2(p)} = \frac{\sigma}{102400 \sqrt{n} \log^2(p)}.$$

This proves the claim when  $p = \max\{3n, \lceil \sqrt{n}/\sigma \rceil\}$ . As in the proof of (10), this is sufficient to prove the claim for any larger value of  $p$ .

□



### 3.2 Proof of Theorem 2

*Theorem 2.* Fix any  $n \geq 30$ ,  $p \geq 3n$ , and  $\sigma \geq 0$ . Then there exists a 2-sparse  $\beta^* \in \mathbb{R}^p$  with  $\frac{1}{2} \leq \|\beta^*\|_2 \leq \|\beta^*\|_1 \leq 1$ , and a positive semi-definite  $\Sigma \in \mathbb{R}^{p \times p}$  with  $\|\Sigma\|_{sp} \leq 2$  and  $\Sigma_{\text{Support}(\beta^*)} = \mathbf{I}_{\text{Support}(\beta^*)}$ , such that for any sample, for all  $B \geq 0$ ,

$$\mathbb{E} \left[ \left( y - x^T \hat{\beta}^B \right)^2 \right] \geq \sigma^2 + \frac{1}{288n \log^2(3n)}. \quad (13)$$

Additionally, if  $100 \leq \sqrt{n}/\sigma \leq p-3$ , then with probability at least  $\frac{1}{2}$  over the sample, for all  $B \geq 0$ ,

$$\mathbb{E} \left[ \left( y - x^T \hat{\beta}^B \right)^2 \right] \geq \sigma^2 + \frac{\sigma}{409600 \sqrt{n} \log^2(\max\{3n, \lceil \sqrt{n}/\sigma \rceil\})}. \quad (14)$$

Here  $\hat{\beta}^B = \arg \min_{\|\beta\|_1 \leq B} \sum_i (y^{(i)} - x^{(i)T} \beta)^2$ , where  $(x^{(i)}, y^{(i)})$  are i.i.d. samples from the multivariate Gaussian distribution defined by drawing  $x^{(i)} \sim N(0, \Sigma)$  and  $y^{(i)} \sim N(x^{(i)T} \beta^*, \sigma^2)$ . The expectations are taken over a new sample  $(x, y)$  drawn from the same distribution, independently of the training set  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$ . (For each  $B \geq 0$ , if  $\hat{\beta}^B$  is not unique, then we show that the inequalities hold for some choice of  $\hat{\beta}^B$ .)

*Proof.* Let  $w_1, w_2, u_1, \dots, u_{p-3} \stackrel{iid}{\sim} N(0, 1)$ . Define

$$\tau = \frac{1}{4 \log p} \cdot \left( \frac{1}{1}, \frac{1}{2}, \dots, \frac{1}{p-3} \right) \in \mathbb{R}^{p-3}.$$

Since  $p \geq 90$ ,  $\|\tau\|_1 \leq \frac{1}{3}$  and  $\|\tau\|_2^2 < \frac{1}{9 \log^2(p)} \leq 0.01$ . Now we define an additional covariate as a linear combination of the others:

$$v = \frac{1}{\sqrt{2}} (w_1 + w_2) \cdot \sqrt{1 - \|\tau\|_2^2} - u^T \tau.$$

Now define  $x = (u_1, \dots, u_{p-3}, v, w_1, w_2)$ . Let  $\Sigma = \text{Cov}(x)$ , and note that  $\sigma_{\max} = \|\Sigma\|_2 \leq 2$ .

Define

$$\beta_{\text{sparse}}^* = \left( \mathbf{0}_{p-3}, 0, \frac{1}{2}, \frac{1}{2} \right), \quad \beta_{\text{dense}}^* = \left( \frac{1}{\sqrt{2(1 - \|\tau\|_2^2)}} \cdot \tau, \frac{1}{\sqrt{2(1 - \|\tau\|_2^2)}}, 0, 0 \right).$$

and

$$y^{(i)} = \frac{1}{2} (w_1 + w_2) = x^{(i)T} \beta_{\text{sparse}}^* = x^{(i)T} \beta_{\text{dense}}^*.$$

Note that  $\beta_{\text{sparse}}^*$  and  $\beta_{\text{dense}}^*$  are both optimal predictors. Since  $\|\beta_{\text{sparse}}^*\|_1 = 1$ ,  $\beta_{\text{sparse}}^{*T} \Sigma \beta_{\text{sparse}}^* = \|\beta_{\text{sparse}}^*\|_2^2 = \frac{1}{2}$ , and  $\beta_{\text{sparse}}^*$  is 2-sparse, this distribution satisfies the desired assumptions. However,

$$\|\beta_{\text{dense}}^*\|_1 = \frac{1}{\sqrt{2(1 - \|\tau\|_2^2)}} (1 + \|\tau\|_1) \approx \frac{4}{3\sqrt{2}} < 1,$$

and so in a sense  $\beta_{\text{dense}}^*$  will be preferred to  $\beta_{\text{sparse}}^*$  in  $\ell_1$ -regularized regression, thus leading to the same arguments as in the proof of Theorem 1.

By Lemma 1, for any  $B \geq 0$ , the set  $\arg \min_{\|\beta\|_1 \leq B} \sum_i (y^{(i)} - x^{(i)T} \beta)^2$  must include a Pareto-optimal vector  $\hat{\beta}^B$  with  $\|\hat{\beta}^B\|_0 \leq n$ . Therefore, it is sufficient to show that bounds (13) and (14) hold for all Pareto-optimal vectors  $\hat{\beta}$  with  $\|\hat{\beta}\|_0 \leq n$ . For each such  $\hat{\beta}$ , we use the notation

$$\hat{\beta} = (\hat{\beta}_u, \hat{\beta}_{w_1}, \hat{\beta}_{w_2}, \hat{\beta}_v) \in \mathbb{R}^{p-3} \times \mathbb{R} \times \mathbb{R} \times \mathbb{R}.$$

Observe that, by definition of the covariates,

$$\begin{aligned} & (\hat{\beta} - \beta_{sparse}^*)^T \Sigma (\hat{\beta} - \beta_{sparse}^*) \\ &= \underbrace{\|\hat{\beta}_u - \tau \hat{\beta}_v\|_2^2}_{\text{(Term 1)}} + \underbrace{\left( \hat{\beta}_{w_1} + \frac{1}{\sqrt{2}} \sqrt{1 - \|\tau\|_2^2} \hat{\beta}_v - \frac{1}{\sqrt{2}} \right)^2}_{\text{(Term 2)}} + \underbrace{\left( \hat{\beta}_{w_2} + \frac{1}{\sqrt{2}} \sqrt{1 - \|\tau\|_2^2} \hat{\beta}_v - \frac{1}{\sqrt{2}} \right)^2}_{\text{(Term 3)}}. \end{aligned} \quad (15)$$

The remainder of the proof is organized as follows. First, we prove bounds (13) and (14) for any Pareto-optimal  $\hat{\beta}$  with  $\hat{\beta}_v \leq \frac{1}{3}$ . Next, we prove the bound (13) for any Pareto-optimal  $\hat{\beta}$  with  $\|\hat{\beta}\|_0 \leq n$  and  $\hat{\beta}_v > \frac{1}{3}$ . Finally, we prove the bound (14) for any Pareto-optimal  $\hat{\beta}$  with  $\|\hat{\beta}\|_0 \leq n$  and  $\hat{\beta}_v > \frac{1}{3}$ .

**Proof of (13) and (14) when  $\hat{\beta}_v \leq \frac{1}{3}$ .** Consider any Pareto-optimal  $\hat{\beta}$  with  $\hat{\beta}_v \leq \frac{1}{3}$ . First, suppose that  $\hat{\beta}_{w_1}, \hat{\beta}_{w_2} \geq \frac{1}{2\sqrt{2}}$ . Let

$$\tilde{\beta} = \left( \hat{\beta}_u + \frac{1}{2\sqrt{1-\|\tau\|_2^2}} \cdot \tau, \hat{\beta}_{w_1} - \frac{1}{2\sqrt{2}}, \hat{\beta}_{w_2} - \frac{1}{2\sqrt{2}}, \hat{\beta}_v + \frac{1}{2\sqrt{1-\|\tau\|_2^2}} \right).$$

By the definition of the covariates,  $x^{(i)T} \hat{\beta} = x^{(i)T} \tilde{\beta}$  for all  $i$ . We will now show that  $\|\tilde{\beta}\|_1 < \|\hat{\beta}\|_1$ . We have

$$\begin{aligned} \|\tilde{\beta}\|_1 &= \|\tilde{\beta}_u\|_1 + |\tilde{\beta}_{w_1}| + |\tilde{\beta}_{w_2}| + |\tilde{\beta}_v| \\ &= \left\| \hat{\beta}_u + \frac{1}{2\sqrt{1-\|\tau\|_2^2}} \cdot \tau \right\|_1 + \left| \hat{\beta}_{w_1} - \frac{1}{2\sqrt{2}} \right| + \left| \hat{\beta}_{w_2} - \frac{1}{2\sqrt{2}} \right| + \left| \hat{\beta}_v + \frac{1}{2\sqrt{1-\|\tau\|_2^2}} \right| \\ &= \left\| \hat{\beta}_u + \frac{1}{2\sqrt{1-\|\tau\|_2^2}} \cdot \tau \right\|_1 + \left| \hat{\beta}_{w_1} \right| + \left| \hat{\beta}_{w_2} \right| - \frac{1}{\sqrt{2}} + \left| \hat{\beta}_v + \frac{1}{2\sqrt{1-\|\tau\|_2^2}} \right| \\ &\leq \left\| \hat{\beta}_u \right\|_1 + \frac{1}{2\sqrt{1-\|\tau\|_2^2}} \cdot \|\tau\|_1 + \left| \hat{\beta}_{w_1} \right| + \left| \hat{\beta}_{w_2} \right| - \frac{1}{\sqrt{2}} + \left| \hat{\beta}_v \right| + \frac{1}{2\sqrt{1-\|\tau\|_2^2}} \\ &= \|\hat{\beta}\|_1 - \frac{1}{\sqrt{2}} + \frac{1}{2\sqrt{1-\|\tau\|_2^2}} + \frac{\|\tau\|_1}{2\sqrt{1-\|\tau\|_2^2}} \leq \|\hat{\beta}\|_1 - \frac{1}{\sqrt{2}} + \frac{1}{2\sqrt{1-0.01^2}} + \frac{0.3}{2\sqrt{1-0.01^2}} \leq \|\hat{\beta}\|_1 - 0.05. \end{aligned}$$

Therefore, this case leads to a contradiction, since we have constructed a coefficient vector  $\tilde{\beta}$  with zero error on the training set, and lower  $\ell_1$ -norm than  $\hat{\beta}$ . Therefore, we must have either  $\hat{\beta}_{w_1} < \frac{1}{2\sqrt{2}}$  or  $\hat{\beta}_{w_2} < \frac{1}{2\sqrt{2}}$ . Without loss of generality, we assume  $\hat{\beta}_{w_1} < \frac{1}{2\sqrt{2}}$ .

Then

$$\hat{\beta}_{w_1} + \frac{1}{\sqrt{2}} \sqrt{1 - \|\tau\|_2^2} \hat{\beta}_v - \frac{1}{\sqrt{2}} < \frac{1}{2\sqrt{2}} + \frac{1}{\sqrt{2}} \cdot \frac{1}{3} - \frac{1}{\sqrt{2}} \leq -\frac{1}{6\sqrt{2}},$$

and so by (Term 2) in (15) above,

$$(\hat{\beta} - \beta_{sparse}^*)^T \Sigma (\hat{\beta} - \beta_{sparse}^*) \geq \left( \hat{\beta}_{w_1} + \frac{1}{\sqrt{2}} \sqrt{1 - \|\tau\|_2^2} \hat{\beta}_v - \frac{1}{\sqrt{2}} \right)^2 \geq \frac{1}{72}.$$

This is sufficient to show that both (13) and (14) are satisfied.

**Proof of (13) when  $\hat{\beta}_v > \frac{1}{3}$ .** Consider any Pareto-optimal  $\hat{\beta}$  with  $\|\hat{\beta}\|_0 \leq n$  and  $\hat{\beta}_v > \frac{1}{3}$ . We then have

$$\begin{aligned} \|\hat{\beta}_u - \tau \hat{\beta}_v\|_2^2 &= \sum_{j=1}^{p-3} \left( \hat{\beta}_{u_j} - \frac{1}{4 \log(p)} \cdot \frac{1}{j} \cdot \hat{\beta}_v \right)^2 \geq \sum_{j \in \{1, \dots, p-3\} : \hat{\beta}_{u_j} = 0} \left( \frac{1}{4 \log(p)} \cdot \frac{1}{j} \cdot \hat{\beta}_v \right)^2 \\ &\geq \frac{\hat{\beta}_v^2}{16 \log^2(p)} \cdot \sum_{j=n}^{p-3} \frac{1}{j^2} \geq \frac{1}{144 \log^2(p)} \cdot \int_{x=n}^{p-3} \frac{1}{x^2} dx = \frac{1}{144 \log^2(p)} \cdot \left( \frac{1}{n} - \frac{1}{p-3} \right) \geq \frac{1}{288n \log^2(p)}. \end{aligned}$$

But, considering (Term 1) in (15) above, this proves that

$$(\hat{\beta} - \beta_{sparse}^*)^T \Sigma (\hat{\beta} - \beta_{sparse}^*) \geq \frac{1}{288n \log^2(p)}.$$

This proves the claim when  $p = 3n$ . As in the proof of Theorem 1, this is sufficient to prove the claim for any larger value of  $p$ .

**Proof of (14) when  $\hat{\beta}_v > \frac{1}{3}$ .** By Lemma 1 of Laurent and Massart [2000], with probability at least  $1 - e^{-0.0625n} \geq 0.75$ ,  $\|z\|_2^2 \sim \chi_n^2 \geq 0.5n$ . For the remainder of the proof, we treat  $z \in \mathbb{R}^n$  as a fixed vector, and assume  $\|z\|_2^2 \geq 0.5n$ .

Assume that (8) holds for all  $J \subset [p]$  with  $|J| = n$ , and (9) holds for all  $J_1 \subset [\lceil \sqrt{n}/\sigma \rceil]$  with  $|J_1| \geq \sqrt{n}/2\sigma$ . (By Lemmas 2 and 3, this is true with probability at least  $1 - 2e^{-n \log(p)} - e^{-0.015\sigma^{-1}\sqrt{n}} \geq 0.75$ .) Consider any Pareto-optimal  $\hat{\beta}$  with  $\|\hat{\beta}\|_0 \leq n$  and  $\hat{\beta}_v > \frac{1}{3}$ . First, suppose that

$$\left| \left\{ j \in [\lceil \sqrt{n}/\sigma \rceil] : \hat{\beta}_j > 0 \right\} \right| < \frac{\sqrt{n}}{2\sigma}.$$

Then

$$\begin{aligned} \|\hat{\beta}_u - \tau \hat{\beta}_v\|_2^2 &= \sum_{j=1}^{p-3} \left( \hat{\beta}_{u_j} - \frac{1}{4 \log(p)} \cdot \frac{1}{j} \cdot \hat{\beta}_v \right)^2 \geq \sum_{j \in [\lceil \sqrt{n}/\sigma \rceil] : \hat{\beta}_{u_j} \leq 0} \left( \frac{1}{4 \log(p)} \cdot \frac{1}{j} \cdot \hat{\beta}_v \right)^2 \\ &= \frac{\hat{\beta}_v^2}{16 \log^2(p)} \sum_{j \in [\lceil \sqrt{n}/\sigma \rceil] : \hat{\beta}_{u_j} \leq 0} \frac{1}{j^2} \geq \frac{1}{144 \log^2(p)} \sum_{j=\lceil \sqrt{n}/2\sigma \rceil}^{\lceil \sqrt{n}/\sigma \rceil} \frac{1}{j^2} \\ &\geq \frac{1}{144 \log^2(p)} \int_{x=\lceil \sqrt{n}/2\sigma \rceil}^{2\lceil \sqrt{n}/2\sigma \rceil} \frac{1}{x^2} dx = \frac{1}{144 \log^2(p)} \left( \frac{1}{\lceil \sqrt{n}/2\sigma \rceil} - \frac{1}{2\lceil \sqrt{n}/2\sigma \rceil} \right) = \frac{1}{144 \log^2(p) \cdot 2\lceil \sqrt{n}/2\sigma \rceil} \geq \frac{\sigma}{288\sqrt{n} \log^2(p)}. \end{aligned}$$

Considering (Term 1) in (15), this proves that

$$(\hat{\beta} - \beta_{sparse}^*)^T \Sigma (\hat{\beta} - \beta_{sparse}^*) \geq \|\hat{\beta}_u - \tau \hat{\beta}_v\|_2^2 \geq \frac{\sigma}{288\sqrt{n} \log^2(p)}.$$

Next, suppose instead that

$$\left| \left\{ j \in [\lceil \sqrt{n}/\sigma \rceil] : \hat{\beta}_j > 0 \right\} \right| \geq \frac{\sqrt{n}}{2\sigma}.$$

Define  $J_1 = \{j \in [\lceil \sqrt{n}/\sigma \rceil] : \hat{\beta}_j > 0\}$ , and fix any  $J \supset \text{Support}(\hat{\beta})$  with  $|J| = n$ . Since  $\hat{\beta}$  is Pareto-optimal with positive entries  $\hat{\beta}_j$  for all  $j \in J_1$ , we have

$$\frac{\partial}{\partial(\beta_{J_1})} \|\hat{\beta}\|_1 = \mathbf{1}_{J_1}.$$

Therefore, by the theory of Lagrange multipliers, we must have  $X_{J_1}^T y - X_{J_1}^T X \hat{\beta} = C \cdot \mathbf{1}_{J_1}$ , for some  $C \in \mathbb{R}$ . We then have

$$X_{J_1}^T X(\hat{\beta} - \beta^*) = \sigma \cdot X_{J_1}^T z - C \cdot \mathbf{1}_{J_1}. \quad (16)$$

By (8), the norm of the left-hand side of (16) can be bounded from above as

$$\left\| X_{J_1}^T X(\hat{\beta} - \beta^*) \right\|_2 \leq \left\| X_{J_1}^T X(\hat{\beta} - \beta^*) \right\|_2 \leq \|\Sigma\|_{sp} \cdot 16\sqrt{2} \cdot n \log(p) \cdot \sqrt{(\hat{\beta} - \beta^*)^T \Sigma (\hat{\beta} - \beta^*)}.$$

By (9), the norm of the right-hand side of (16) can be bounded from below as

$$\left\| \sigma \cdot X_{J_1}^T z - C \cdot \mathbf{1}_{J_1} \right\|_2 \geq \sigma \cdot \left\| \text{Proj}_{\mathbf{1}_{J_1}^\perp} X_{J_1}^T z \right\|_2 \geq \sigma \sqrt{\frac{\lambda_{\min}^2(\Sigma_{J_1}) n^{3/2}}{200\sigma}}.$$

Therefore, returning to (16), we have

$$\begin{aligned} \|\Sigma\|_{sp} \cdot 16\sqrt{2} \cdot n \log(p) \cdot \sqrt{(\hat{\beta} - \beta^*)^T \Sigma (\hat{\beta} - \beta^*)} &\geq \left\| X_{J_1}^T X(\hat{\beta} - \beta^*) \right\|_2 \\ &= \left\| \sigma \cdot X_{J_1}^T z - C \cdot \mathbf{1}_{J_1} \right\|_2 \geq \sigma \sqrt{\frac{\lambda_{\min}^2(\Sigma_{J_1}) n^{3/2}}{200\sigma}}. \end{aligned}$$

Therefore,

$$(\hat{\beta} - \beta^*)^T \Sigma (\hat{\beta} - \beta^*) \geq \frac{\sigma \cdot \lambda_{\min}^2(\Sigma_{J_1})}{102400 \|\Sigma\|_{sp}^2 \cdot \sqrt{n} \log^2(p)} = \frac{\sigma}{409600 \sqrt{n} \log^2(p)}.$$

This proves the claim when  $p = \max\{3n, \lceil \sqrt{n}/\sigma \rceil\}$ . As in the proof of Theorem 1, this is sufficient to prove the claim for any larger value of  $p$ .

□

## 4 Proofs for Lemmas

### 4.1 Proof of Lemma 2

Fix any  $J \subset [p]$  with  $|J| = n$ . We will show that, with probability at least  $1 - 2e^{-2n \log(p)}$ ,

$$\left\| X_J^T X(\tilde{\beta} - \beta^*) \right\|_2 \leq \|\Sigma\|_{sp} \cdot 16\sqrt{2} \cdot n \log(p) \cdot \sqrt{(\tilde{\beta} - \beta^*)^T \Sigma (\tilde{\beta} - \beta^*)} \text{ for all } \tilde{\beta} \in \mathbb{R}^p \text{ with } \tilde{\beta}_{\bar{J}} = 0.$$

Since there are  $\binom{p}{n} \leq p^n$  choices for the set  $J$ , this will be sufficient to prove the lemma.

Reorder the covariates to write  $\Sigma = \begin{pmatrix} \Sigma_{JJ} & \Sigma_{J\bar{J}} \\ \Sigma_{\bar{J}J} & \Sigma_{\bar{J}\bar{J}} \end{pmatrix}$ . Choose a Cholesky decomposition

$$\Sigma = \begin{pmatrix} U & V \\ \mathbf{0} & W \end{pmatrix}^T \begin{pmatrix} U & V \\ \mathbf{0} & W \end{pmatrix}.$$

Let  $a^{(i)} \sim N(0, \mathbf{I}_p)$ . Then  $\begin{pmatrix} U & V \\ \mathbf{0} & W \end{pmatrix}^T a^{(i)} \sim N(0, \Sigma)$ , and so

$$\begin{pmatrix} X_J & X_{\bar{J}} \end{pmatrix} \stackrel{\mathcal{D}}{=} \begin{pmatrix} A_J U & A_J V + A_{\bar{J}} W \end{pmatrix},$$

where the matrix  $A$  has entries  $A_{ij} = a_j^{(i)}$ , and  $A_J$  consists of the columns of  $A$  indexed by  $j \in J$ . We then have

$$\begin{aligned} X_J^T X (\tilde{\beta} - \beta^*) &\stackrel{\mathcal{D}}{=} U^T A_J^T \left( A_J U (\tilde{\beta} - \beta^*)_J + (A_J V + A_{\bar{J}} W) (\tilde{\beta} - \beta^*)_{\bar{J}} \right) \\ &= U^T A_J^T \left( A_J U (\tilde{\beta} - \beta^*)_J - (A_J V + A_{\bar{J}} W) \beta_{\bar{J}}^* \right) = U^T A_J^T \left( A_J \left( U (\tilde{\beta} - \beta^*)_J - V \beta_{\bar{J}}^* \right) - A_{\bar{J}} W \beta_{\bar{J}}^* \right) \end{aligned}$$

Below, we will show that

$$\|A_J\|_{sp} \leq \sqrt{16n \log(p)} \text{ with probability at least } 1 - e^{-2n \log(p)}, \quad (17)$$

$$\text{and } \|A_{\bar{J}} W \beta_{\bar{J}}^*\|_2 \leq \sqrt{16n \log(p)} \cdot \|W \beta_{\bar{J}}^*\|_2 \text{ with probability at least } 1 - e^{-2n \log(p)}. \quad (18)$$

Assuming that these bounds hold. Then for any  $\tilde{\beta} \in \mathbb{R}^p$  with  $\tilde{\beta}_{\bar{J}} = 0$ , we have

$$\begin{aligned} &\left\| U^T A_J^T \left( A_J \left( U (\tilde{\beta} - \beta^*)_J - V \beta_{\bar{J}}^* \right) - A_{\bar{J}} W \beta_{\bar{J}}^* \right) \right\|_2 \\ &\leq \|U\|_{sp} \cdot \|A_J\|_{sp} \cdot \left( \|A_J\|_{sp} \cdot \left\| U (\tilde{\beta} - \beta^*)_J - V \beta_{\bar{J}}^* \right\|_2 + \|A_{\bar{J}} W \beta_{\bar{J}}^*\|_2 \right) \\ &\leq \|\Sigma\|_{sp} \cdot \sqrt{16n \log(p)} \cdot \left( \sqrt{16n \log(p)} \cdot \left\| U (\tilde{\beta} - \beta^*)_J - V \beta_{\bar{J}}^* \right\|_2 + \sqrt{16n \log(p)} \cdot \|W \beta_{\bar{J}}^*\|_2 \right) \\ &= \|\Sigma\|_{sp} \cdot 16n \log(p) \cdot \left( \left\| U (\tilde{\beta} - \beta^*)_J - V \beta_{\bar{J}}^* \right\|_2 + \|W \beta_{\bar{J}}^*\|_2 \right) \\ &\leq \|\Sigma\|_{sp} \cdot 16n \log(p) \cdot \sqrt{2} \cdot \left( \left\| U (\tilde{\beta} - \beta^*)_J - V \beta_{\bar{J}}^* \right\|_2^2 + \|W \beta_{\bar{J}}^*\|_2^2 \right)^{1/2} \\ &= \|\Sigma\|_{sp} \cdot 16\sqrt{2} \cdot n \log(p) \cdot \sqrt{(\tilde{\beta} - \beta^*)^T \Sigma (\tilde{\beta} - \beta^*)}. \end{aligned}$$

We conclude by proving (17) and (18). We first prove (17) using a construction from Keshavan et al. [2010]. First, define  $\mathcal{U} = \left\{ u \in \left( \frac{1}{8\sqrt{n}} \mathbb{Z} \right)^n : \|u\|_2 \leq 1 \right\}$ . By Remark 5.1 in Keshavan et al. [2010],

$$\|A_J\|_{sp} \leq \sqrt{2} \sup_{u, v \in \mathcal{U}} |u^T A_J v|.$$

For any  $u, v \in \mathcal{U}$ ,

$$u^T A_J v = \sum_{i=1}^n \sum_{j \in J} u_i v_j A_{ij} \sim N(0, \|u\|_2^2 \|v\|_2^2),$$

therefore,

$$\Pr \left( |u^T A_J v| \geq \sqrt{8n \log(p)} \right) \leq \Pr \left( |N(0, 1)| \geq \sqrt{8n \log(p)} \right) \leq e^{-4n \log(p)}.$$

Furthermore,  $|\mathcal{U}| \leq (2 \lceil 8\sqrt{n} \rceil + 1)^n \leq p^n$ . So,

$$\begin{aligned} \Pr \left( \|A_J\|_{sp} \leq \sqrt{16n \log(p)} \right) &\leq \Pr \left( |u^T A_J v| \leq \sqrt{8 \cdot n \log(p)} \text{ for all } u, v \in \mathcal{U} \right) \\ &\geq 1 - p^{2n} e^{-4n \log(p)} \geq 1 - e^{-2n \log(p)}. \end{aligned}$$

Next we prove (18). We have

$$\|A_{\mathcal{J}} W \beta_{\mathcal{J}}^*\|_2^2 = \sum_i \left( \sum_{j \in \mathcal{J}} a_j^{(i)} (W \beta_{\mathcal{J}}^*)_j \right)^2 \stackrel{\mathcal{D}}{=} \|W \beta_{\mathcal{J}}^*\|_2^2 \cdot \chi_n^2.$$

By Lemma 1 of Laurent and Massart [2000],  $\Pr(\chi_n^2 \geq 16n \log(p)) \leq e^{-2n \log(p)}$ . This is sufficient.

## 4.2 Proof of Lemma 3

Choose any  $J_2 \subset J_1$  with  $|J_2| = \lceil \sqrt{n}/2\sigma \rceil$ . Observe that  $\left\| \text{Proj}_{\mathbf{1}_{J_1}}^\perp (X_{J_1}^T z) \right\|_2^2 \geq \left\| \text{Proj}_{\mathbf{1}_{J_2}}^\perp (X_{J_2}^T z) \right\|_2^2$ , and so it is sufficient to only consider the sets  $J_2$  of size  $\lceil \sqrt{n}/2\sigma \rceil$ .

Fix any  $J_2 \subset [\lceil \sqrt{n}/\sigma \rceil]$  with  $|J_2| = \lceil \sqrt{n}/2\sigma \rceil$ . Let  $P \in \mathbb{R}^{\lceil \sqrt{n}/2\sigma \rceil \times \lceil \sqrt{n}/2\sigma \rceil}$  be the orthogonal projection matrix corresponding to  $\text{Proj}_{\mathbf{1}_{J_2}}^\perp(\cdot)$ . Write  $P \Sigma_{J_2} P = A A^T$  for  $A \in \mathbb{R}^{\lceil \sqrt{n}/2\sigma \rceil \times (\lceil \sqrt{n}/2\sigma \rceil - 1)}$ . Then  $(X_{J_2}^T z) \sim N(0, \|z\|_2^2 \cdot \Sigma_{J_2})$  and so  $\text{Proj}_{\mathbf{1}_{J_2}}^\perp (X_{J_2}^T z) \sim N(0, \|z\|_2^2 \cdot P \Sigma_{J_2} P)$ , and therefore  $\text{Proj}_{\mathbf{1}_{J_2}}^\perp (X_{J_2}^T z) \stackrel{\mathcal{D}}{=} \|z\|_2^2 \cdot A u$  for  $u \sim N(0, \mathbf{I}_{\lceil \sqrt{n}/2\sigma \rceil - 1})$ . By examining the definition of  $A$ , we see that  $u^T (A^T A) u \geq \|u\|_2^2 \cdot \lambda_{\min}^2(\Sigma_{\lceil \sqrt{n}/\sigma \rceil})$ , therefore,

$$\left\| \text{Proj}_{\mathbf{1}_{J_2}}^\perp (X_{J_2}^T z) \right\|_2^2 \stackrel{\mathcal{D}}{=} \|z\|_2^2 \cdot \|A u\|_2^2 \stackrel{\mathcal{D}}{\geq} 0.5n \cdot \lambda_{\min}^2(\Sigma_{\lceil \sqrt{n}/\sigma \rceil}) \cdot \chi_{\lceil \sqrt{n}/2\sigma \rceil - 1}^2.$$

Furthermore, the number of such sets  $J_2$  is bounded by  $2^{\lceil \sqrt{n}/\sigma \rceil}$ . By the chi-square tail bounds from Foygel and Drton [2010], using the assumption that  $\sqrt{n}/\sigma \geq 100$ , we have

$$\begin{aligned} \Pr \left( \chi_{\lceil \sqrt{n}/2\sigma \rceil - 1}^2 \leq \frac{\sqrt{n}}{100\sigma} \right) &\leq \Pr \left( \chi_{\lceil \sqrt{n}/2\sigma \rceil - 1}^2 \leq 0.02 \cdot \frac{50}{49} \cdot (\lceil \sqrt{n}/2\sigma \rceil - 1) \right) \\ &\leq \exp \left\{ \frac{1}{2} (\lceil \sqrt{n}/2\sigma \rceil - 2) \left( 1 - 0.02 \cdot \frac{50}{49} + \log \left( 0.02 \cdot \frac{50}{49} \right) \right) \right\} \leq \exp \left\{ \frac{1}{2} (\lceil \sqrt{n}/2\sigma \rceil \cdot \frac{48}{50}) \left( 1 - 0.02 \cdot \frac{50}{49} + \log \left( 0.02 \cdot \frac{50}{49} \right) \right) \right\} \\ &\leq e^{-0.7084 \lceil \sqrt{n}/\sigma \rceil}. \end{aligned}$$

Therefore,

$$\begin{aligned} \Pr \left( \exists J_1 \subset [\lceil \sqrt{n} \rceil], |J_1| \geq \frac{\sqrt{n}}{2}, \left\| \text{Proj}_{\mathbf{1}_{J_1}}^\perp (X_{J_1}^T z) \right\|_2^2 \leq n \lambda_{\min}^2(\Sigma_{\lceil \sqrt{n} \rceil}) \cdot \frac{\sqrt{n}}{200\sigma} \right) \\ \leq \Pr \left( \exists J_2 \subset [\lceil \sqrt{n} \rceil], |J_2| = \frac{\lceil \sqrt{n} \rceil}{2}, \left\| \text{Proj}_{\mathbf{1}_{J_2}}^\perp (X_{J_2}^T z) \right\|_2^2 \leq n \lambda_{\min}^2(\Sigma_{\lceil \sqrt{n} \rceil}) \cdot \frac{\sqrt{n}}{200\sigma} \right) \\ \leq 2^{\lceil \sqrt{n}/\sigma \rceil} \cdot \Pr \left( \chi_{\lceil \sqrt{n}/2\sigma \rceil - 1}^2 \leq \frac{\sqrt{n}}{100\sigma} \right) \leq 2^{\lceil \sqrt{n}/\sigma \rceil} \cdot e^{-0.7084 \lceil \sqrt{n}/\sigma \rceil} \leq e^{-0.015 \lceil \sqrt{n}/\sigma \rceil} \leq e^{-0.015 \sigma^{-1} \sqrt{n}}. \end{aligned}$$

## References

P.J. Bickel, Y. Ritov, and A.B. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.

- T.T. Cai, G. Xu, and J. Zhang. On recovery of sparse signals via  $\ell_1$  minimization. *IEEE Transactions on Information Theory*, 55(7):3388–3397, 2009.
- E.J. Candes and T. Tao. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203–4215, 2005. ISSN 0018-9448.
- R. Foygel and M. Drton. Extended Bayesian information criteria for Gaussian graphical models. *Advances in Neural Information Processing Systems*, 23:604–612, 2010.
- R.H. Keshavan, A. Montanari, and S. Oh. Matrix completion from noisy entries. *Journal of Machine Learning Research*, 11:2057–2078, 2010.
- V. Koltchinskii. The Dantzig selector and sparsity oracle inequalities. *Bernoulli*, 15(3):799–828, 2009.
- B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, 28(5):1302–1338, 2000.
- S. Rosset, J. Zhu, and T. Hastie. Boosting as a regularized path to a maximum margin classifier. *Journal of Machine Learning Research*, 5:941–973, 2004.
- N. Srebro, K. Sridharan, and A. Tewari. Smoothness, low noise and fast rates. *Advances in Neural Information Processing Systems*, 23:2199–2207, 2010.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996.